

Compilerconstructie

najaar 2012

<http://www.liacs.nl/home/rvv11let/coco/>

Rudy van Vliet

Kamer 124 Snellius, tel. 071-527 5777
rvv1iet(at)liacs.nl

college 2, dinsdag 11 september 2012

Lexical Analysis

1

2.6 Lexical Analyser

Reads and converts the input into a stream of tokens to be analysed by the parser

Lexeme: Sequence of input characters comprising single token

Typical tasks of the lexical analyser

- Remove white space and comments
- Encode constants as tokens:
31 + 28 + 59 → (num, 31) (+) (num, 28) (+) (num, 59)
- Recognize keywords
- Recognize identifiers:
count = count + increment; →
(id, "count") (=) (id "count") (+) (id, "increment") (:)

Lexical analyser may need to read ahead (with input buffer)

2

2.7 Symbol Table

- Symbol table holds information about *source-program constructs* (e.g., identifiers)
 - string
 - additional information (type, position in storage)
- Symbol table is globally accessible (to all phases of compiler)
- Information is collected incrementally by analysis phases, and used by synthesis phases
- Implementation by Hashtable, with methods
 - *put (String, Symbol)*
 - *get (String)*

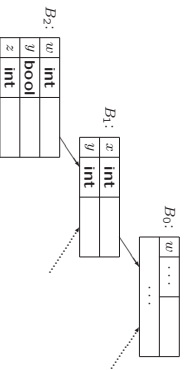
3

Symbol Table Per Scope

The same identifier may be declared more than once

```

1) { int x; int y;
2) { int w; bool y; int z;
3) ... w ...; ... x ...; ... y ...; ... z ...;
4) }
5) ... w ...; ... x ...; ... y ...;
6) }
    
```



5

The Use of Symbol Tables

```

program → { top = null; }
        block
block → '{' { saved = top;
           decls stmts '}' { top = new Env(top);
                           top = saved;
decls → decls decl
           | ε
decl → type id: { s = new Symbol(
                s.type = type.lexeme,
                top.put(id.lexeme, s);
                }
    
```

In book (edition 2) extended for real translation

7

Symbol Table Per Scope

The same identifier may be declared more than once

```

1) { int x; int y;
2) { int w; bool y; int z;
3) ... w ...; ... x ...; ... y ...; ... z ...;
4) }
5) ... w ...; ... x ...; ... y ...;
6) }
    
```

4

Translation Scheme (Example)

(from college 1)

```

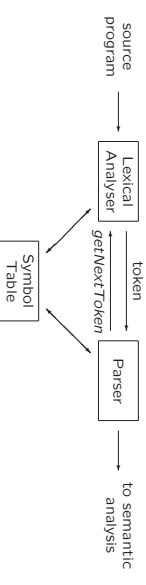
expr → expr1 + term {print('+')}
expr → expr1 - term {print('-')}
expr → term
term → 0 {print('0')}
term → 1 {print('1')}
...
term → 9 {print('9')}
    
```

Example: parse tree for 9 – 5 + 2

Implementation requires postorder traversal

6

3.1 Lexical Analyser - Parser Interaction



8

Lexical Analyser

Reasons why it is a separate phase of a compiler

- Simplifies the design of the compiler
- Provides efficient implementation
 - Systematic techniques to implement lexical analysers (by hand or automatically)
 - Stream buffering methods to scan input
- Improves portability
 - Non-standard symbols and alternate character encodings can be more easily translated

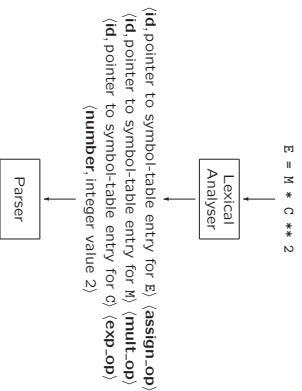
9

Tokens, Patterns and Lexemes

- **Token**: pair of token name and optional attribute value, e.g., `<id, 1>`, `<num, 31>`, `<assign, op>`
- **Lexeme**: specific sequence of characters that makes up token, e.g., `count`, `31`, `=`
- **Pattern**: description of form that lexemes of a token may take, e.g.,
 - if: if
 - comparison**: `< or >` or `<= or >=` or `== or !=`
 - id**: letter followed by letters and digits

10

Attributes for Tokens



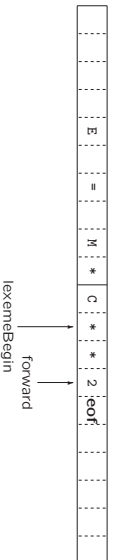
11

3.2 Input Buffering

Use two buffers of size N for input

- Saves time
- Allows for looking ahead one or more characters, e.g., for
 - identifiers: `ifoundit`
 - relational operators: `<=`

Take longest prefix of input that matches any pattern



13

3.3 Specification of Tokens

Regular expressions to specify patterns for tokens

Terminology (from F11)

- An **alphabet** Σ is a finite set of symbols (characters), e.g., `{0,1}`, ASCII, Unicode
- A **string** s is a finite sequence of symbols from Σ
 - $|s|$ denotes the length of string s , e.g., `|banana| = 6`
 - ϵ denotes an empty string: $|\epsilon| = 0$
- A **language** is a set of strings over some fixed alphabet Σ

15

Lexical Errors

- Hard to detect by lexical analyser alone, e.g.,
 - `f1 (a == f(x)) ...`
- What if none of the patterns matches?
 - ‘Panic mode’ recovery: delete characters until you find well-formed token
 - * Delete one character from remaining input
 - * Insert missing character into remaining input
 - * Replace character by another character
 - * Transpose two adjacent characters

12

Implement a Lexical Analyser

- By hand, using transition diagram to specify lexemes
- With a lexical-analyser generator (Lex), using regular expressions to specify lexemes:
 - Regular expressions \rightarrow (non-deterministic) finite automaton \rightarrow deterministic finite automaton
 - Input to ‘driver’

14

String operations

- **Concatenation** of strings x and y is denoted xy e.g., if $x = \text{dog}$ and $y = \text{house}$ then $xy = \text{doghouse}$ $sc = \text{cs} = s$
- **Exponentiation**
 - Define
 - $s^0 = \epsilon$
 - $s^i = s^{i-1}s$ if $i > 0$
 - Then
 - $s^1 = s$
 - $s^2 = ss$
 - $s^3 = sss$

16

Language Operations

- **Union** $L \cup D = \{s \mid s \in L \text{ or } s \in D\}$
- **Concatenation** $LD = \{xy \mid x \in L \text{ and } y \in D\}$
- **Exponentiation** $L^0 = \{\epsilon\}; \quad L^i = L^{i-1}L \quad \text{if } i > 0$
- **Kleene closure** $L^* = \bigcup_{i=0}^{\infty} L^i$
(zero or more concatenation)
- **Positive closure** $L^+ = \bigcup_{i=1}^{\infty} L^i$
(one or more concatenation)

17

Language Operations (Example)

Let alphabets $L = \{A, B, \dots, Z, a, b, \dots, z\}$ and $D = \{0, 1, \dots, 9\}$

- $L \cup D$ is set of letters and digits
- LD is set of strings consisting of a letter followed by a digit
- L^4 is set of all four-letter strings
- L^* is set of all finite strings of letters, including ϵ
- $L(L \cup D)^*$ is set of all strings of letters and digits beginning with a letter ('identifiers')
- D^+ is set of all strings of one or more digits ('nonnegative integers')

18

Regular Expressions (Definition)

- Each regular expression r denotes a language $L(r)$
- Defining rules:
 - ϵ is regular expression, and $L(\epsilon) = \{\epsilon\}$
 - if $a \in \Sigma$, then a is regular expression, and $L(a) = \{a\}$.
 - if r and s are regular expressions, then
 - * $(r) \mid (s)$ is regular expression denoting $L(r) \cup L(s)$
 - * $(r)(s)$ is regular expression denoting $L(r)L(s)$
 - * $(r)^*$ is regular expression denoting $L(r)^*$
 - * $(r)^+$ is regular expression denoting $L(r)^+$

19

20

Regular Expressions (Example)

In C, an identifier is a letter followed by zero or more letters or digits (underscore is considered letter):

$letter_ (letter_ | digit)^*$

Regular Expressions (Example)

- Remove unnecessary parentheses by assuming precedence relation between $*$, concatenation, and $|$, e.g.,
 $(a) \mid ((b)^*(c))$ is equivalent to $a \mid b^*c$
- Let $\Sigma = \{a, b\}$. Then the regular expression:
 - $a \mid b$ denotes the set $\{a, b\}$
 - $(a \mid b)(a \mid b)$ denotes the set $\{aa, ab, ba, bb\}$
 - a^* denotes the set $\{\epsilon, a, aa, aaa, \dots\}$
 - $(a \mid b)^*$ denotes the sets of all strings over $\{a, b\}$
 - $a \mid a^*b$ denotes the string a and all strings consisting of zero or more a 's followed by one b
- If r and s denote the same language L , then $r = s$, e.g., $(a \mid b) = (b \mid a)$

21

22

Regular Definitions

- A **regular definition** is a sequence of definitions of the form:

$$\begin{aligned} d_1 &\rightarrow r_1 \\ d_2 &\rightarrow r_2 \\ &\dots \\ d_n &\rightarrow r_n \end{aligned}$$

where r_i is a regular expression over $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$

- Obtain regular expression over Σ by successively substituting d_j ($j = 1, 2, \dots, n - 1$) in r_{j+1}, \dots, r_n by (r_j)

23

Regular Definitions

- A **regular definition** is a sequence of definitions of the form:
 - $d_1 \rightarrow r_1$
 - $d_2 \rightarrow r_2$
 - \dots
 - $d_n \rightarrow r_n$
 where r_i is a regular expression over $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$
 - Obtain regular expression over Σ by ...

22

Regular Definitions (Example)

- Identifiers in C

$$\begin{aligned} letter_ &\rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z \mid - \\ digit &\rightarrow 0 \mid 1 \mid \dots \mid 9 \\ Id &\rightarrow letter_ (letter_ \mid digit)^* \end{aligned}$$

- Recursion is not allowed

$digit \rightarrow digit(digit)^*$	not OK
$digits \rightarrow digit(digit)^*$	OK

24

Notational Shorthands

- We often use the following shorthands:
 - one-or-more instance of: r^+ = rr^*
 - zero-or-one instance of: $r^?$ = $r | \epsilon$
 - Character classes:
 - $[abd]$ = $a | b | d$
 - $[a-z]$ = $a | b | \dots | z$
- Example, unsigned numbers:
 - 5280, 0.01234, 6.336E4, 1.89E-4

```
digit → [0-9]
digits → digit+
number → digits(.digits)?([E|+|-]?digits)?
```

25

Regular Definitions for Tokens

Regular definitions describing patterns for these tokens

```
digit → [0-9]
digits → digit+
number → digits(.digits)?([E|+|-]?digits)?
letter → [A-Za-z]
id → letter(letter | digit)*
if → if
then → then
else → else
relop → < | > | <= | <= | >= | = | <>
```

Regular definition for white space

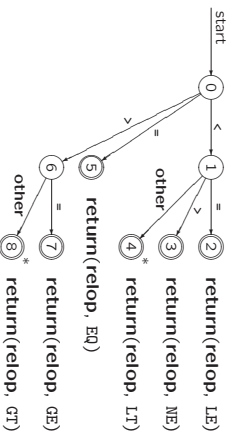
```
ws → (blank | tab | newline)+
```

27

Transition Diagrams

(‘Almost finite automata’)

```
relop → < | > | <= | >= | = | <>
```



Retract input one position, if necessary (*)

29

Transition Diagrams



How to distinguish between identifiers and (reserved) keywords?
Two possibilities:

- Install reserved words in symbol table initially
Used in above diagram
- Separate transition diagram for each keyword
Try these first, before the diagram for identifiers

31

3.4 Recognition of Tokens

Grammar for branching statements:

```
stmt → if expr then stmt
      | if expr then stmt else stmt
      | ε
expr → term relop term
      | term
term → id
      | number
```

Terminals are **if**, **then**, **else**, **relop**, **id** and **number**.
These are the names of the tokens.

26

Lexemes and Their Tokens

Goal:

Lexemes	Token name	Attribute value
Any ws	-	-
if	if	-
then	then	-
else	else	-
Any id	id	pointer to table entry
Any number	number	pointer to table entry
<	relop	LT
<=	relop	LE
=	relop	EQ
>	relop	ME
>=	relop	GT

28

Transition Diagrams

Identifiers and keywords

```
id → letter(letter | digit)*
```



How to distinguish between identifiers and (reserved) keywords?

30

From Diagram to Lexical Analyser

```
TOKEN getRelop ()
{
    TOKEN reToken = new (RELOP);
    while (1)
        /* repeat character processing until a return
        or failure occurs */
        switch (state)
        {
            case 0: c = nextChar();
                    if (c == '<' ) state = 1;
                    else if (c == '>' ) state = 5;
                    else if (c == '<=' ) state = 6;
                    else fail(); /* lexeme is not a relop */
                    break;
            case 1: ...
            ...
            case 8: retract();
                    reToken.attribute = GT;
                    return(reToken);
        }
    }
}
```

32

Entire Lexical Analyser

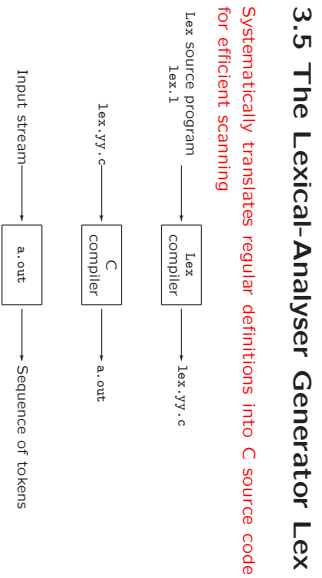
Based on transition diagrams for different tokens
Three possibilities:

- Try transition diagrams sequentially (in right order)
- Run transition diagrams in parallel
Make sure to take longest prefix of input that matches any pattern
- Combine all transition diagrams into one

33

Entire Lexical Analyser

Based on transition diagrams for different tokens
How?



35

Operation of Lexical Analyser

The lexical analyser generated by Lex

- Activated by parser
- Reads input character by character
- Executes action A_i corresponding to pattern P_i
- Typically, A_i returns to the parser
- If not (e.g., in case of white space), proceed to find additional lexemes
- Lexical analyser returns single value: the token name
- Attribute value passed through global variable `yy1val`

37

Lexemes and Their Tokens

Goal:

Lexemes	Token name	Attribute value
Any ws	—	—
if	if	—
then	then	—
else	else	—
Any id	id	pointer to table entry
Any number	number	pointer to table entry
<	relop	LT
<=	relop	LE
=	relop	EQ
>	relop	NE
>=	relop	GT
	relop	GE

39

Entire Lexical Analyser

Based on transition diagrams for different tokens
Three possibilities:

- Try transition diagrams sequentially (in right order)
 - Run transition diagrams in parallel
Make sure to take longest prefix of input that matches any pattern
 - Combine all transition diagrams into one
- A Lex program has the following form
- ```

declarations
%%
translation rules
%%
user defined auxiliary functions

• Translation rules are of the form

Pattern { Action }

Patterns are Lex regular expressions

```

36

## Regular Definitions for Tokens

Regular definitions describing patterns for these tokens

```

digit → [0-9]
digits → digit+
number → digits{?E[+-]?digits}?
letter → [A-Za-z]
id → letter{letter | digit}*
then → then
else → else
relop → < | > | <= | >= | = | <>

```

Regular definition for white space

```
ws → (Blank | tab | newline)+
```

38

## The Lex Program (program.l)

```

/* declarations section */
%{
/* definitions of constants */
#define LT 256
/* macros for LE, EQ, NE, GT, GE,
IF, THEN, ELSE, ID, NUMBER, HELD */
%}

/* regular definitions)
delim [\v\n]
ws {delim}+
letter [A-Za-z]
digit [0-9]
id {letter}{digit}*
number {digit}+(\.){digit}+{E[+-]}{digit}+)?

```

40

## The Lex Program (program.l)

```

%%
/* translation rules section */
{ws} { /* no action and no return */ }
if {return(IF);}
then {return(THEN);}
else {return(ELSE);}
{fd} {yy1val = (int) installNum(); return(NUMBER);}
faunber} {yy1val = (int) installNum(); return(NUMBER);}
"<" {yy1val = LT; return(RELOP);}
"<=" {yy1val = LE; return(RELOP);}
"=" {yy1val = EQ; return(RELOP);}
"<=" {yy1val = NE; return(RELOP);}
">" {yy1val = GT; return(RELOP);}
">=" {yy1val = GE; return(RELOP);}

%%
/* auxiliary functions section */
int installID() {...}
int installNum() {...}

```

41

## Regular expressions in Lex

Operator characters: \ " . ^ \$ [ ] \* + ? { } | /

| Expression | Matches                               | Example |
|------------|---------------------------------------|---------|
| c          | non-operator character                | a       |
| \c         | operator character                    | **      |
| "s"        | string s literally                    | a.*b    |
| .          | any character but newline             | abc     |
| \$         | beginning of a line                   | abc\$   |
| [s]        | any one of the characters in string s | [abc]   |
| [^s]       | any one character not in string s     | [^abc]  |
| [a-c2]     | any one character between c1 and c2   | [a-z]   |
| r*         | zero or more strings matching r       | aa      |
| r+         | one or more strings matching r        | a+      |
| r?         | zero or one string matching r         | a?      |
| r{m, n}    | between m and n occurrences of r      | at{1,5} |
| r1  r2     | an r1 followed by an r2               | ab      |
| r1 / r2    | an r1 or an r2                        | ab b    |
| (r)        | same as r                             | (a b)   |
| /d/        | r1 when followed by r2                | abc/123 |
| {d}        | regular expression defined by d       | {1d}    |

42

## Lex Details

- installID()
  - function to install the lexeme into the symbol table
  - returns pointer to symbol table entry
- yytext
  - pointer to the first character of the lexeme
- yyLeng
  - length of the lexeme
- installNum()
  - similar to installID, but puts numerical constants into a separate table

43

## Compiler constructie

College 2  
Lexical Analysis

Chapters for reading: 2.6, 2.7, 3.1–3.5

45

## Lex Details

- Example: input "\t\tif "
  - Longest initial prefix: "\t\t" = WS
  - No action, so yytext points to 'i' and continue
  - Next lexeme is "if"
  - Token **if** is returned, yytext points to 'i' and yyLeng=2
- Ambiguity and longest pattern matching:
  - Patterns **if** and **{fd}** match lexeme "if"
  - If input is "<=", then lexeme is "<="
- lex program.1
  - gcc lex.yy.c -l1
  - ./a.out < input

44